

MACHINE LEARNING IN BIOINFORMATICS

PROBABILITY BASICS

Philipp Benner

philipp.benner@bam.de

VP.1 - eScience

Federal Institute of Materials Research and Testing (BAM)

April 25, 2024

INTRODUCTION TO PROBABILITY THEORY

ROULETTE WHEEL



ROULETTE WHEEL

Assume we have a fair roulette wheel with 37 segments, of which

- 18 are black
- 18 are red
- 1 is green and labeled with a zero

The red and black segments are labeled with numbers ranging from 1 to 36, where

	even	odd
red	8	10
black	10	8

ROULETTE WHEEL - SIMPLE PROBABILITIES

- What is the probability of black?

$$\text{pr}(\text{black}) = \frac{\#\text{black segments}}{\#\text{segments}} = \frac{18}{37}$$

ROULETTE WHEEL - SIMPLE PROBABILITIES

- What is the probability of black?

$$\text{pr}(\text{black}) = \frac{\#\text{black segments}}{\#\text{segments}} = \frac{18}{37}$$

- What is the probability of black or green?

$$\begin{aligned}\text{pr}(\text{black or green}) &= \frac{\#\text{black and green segments}}{\#\text{segments}} \\ &= \frac{\#\text{black segments}}{\#\text{segments}} + \frac{\#\text{green segments}}{\#\text{segments}} \\ &= \text{pr}(\text{black}) + \text{pr}(\text{green})\end{aligned}$$

This property is called **additivity**

- What is the probability of observing first black and afterwards red?

$$\text{pr}(\text{first black and then red}) = \text{pr}(\text{black})\text{pr}(\text{red})$$

- What is the probability of observing first black and afterwards red?

$$\text{pr}(\text{first black and then red}) = \text{pr}(\text{black})\text{pr}(\text{red})$$

- **and** \Rightarrow "multiplication"
- **or** \Rightarrow "addition"

ROULETTE WHEEL - SIMPLE PROBABILITIES

- What is the probability of a black segment with an even number?

$$\begin{aligned} & \text{pr}(\text{black and even number}) \\ &= \text{pr}(\text{black})\text{pr}(\text{even number}) \\ &= \frac{\# \text{black segments}}{\# \text{segments}} \frac{\# \text{even segments}}{\# \text{segments}} \end{aligned}$$

ROULETTE WHEEL - SIMPLE PROBABILITIES

- What is the probability of a black segment with an even number?

$$\begin{aligned} & \text{pr}(\text{black and even number}) \\ &= \text{pr}(\text{black})\text{pr}(\text{even number}) \\ &= \frac{\# \text{black segments}}{\# \text{segments}} \frac{\# \text{even segments}}{\# \text{segments}} \end{aligned}$$

Wrong!

- Both events are not **independent!** Some black segments are even.

Sample space

The set of all possible outcomes is called **sample space** and typically denoted Ω . The elements of the sample space are called **outcomes** or *samples*

- For our roulette wheel, if we only care about the color of segments, then the sample space is

$$\Omega = \{\text{red, black, green}\}$$

- If we consider both colors and numbers, then

$$\Omega = \{0 : \text{green}, 1 : \text{red}, 2 : \text{black}, \dots\}$$

- Colors and numbers are not independent:

$$\Omega \neq \{\text{red, black, green}\} \times \{0, 1, 2, \dots, 36\}$$

Events

An event E is any subset of Ω , denoted $E \subseteq \Omega$

- We assign probabilities to events $E \subseteq \Omega$
- The probability of "black or green" is denoted

$$\text{pr}(\{\text{black}, \text{green}\})$$

- More formally, we may write $\text{pr}(E)$ for some $E \subseteq \Omega$

PROBABILITY AXIOMS – AXIOM I

- What is the lowest possible probability?
- Assume that

$$\Omega = \{\text{yellow, red, black, green}\}$$

then $\text{pr}(\{\text{yellow}\}) = 0$, since there is no yellow segment

- We could also write $\text{pr}(\emptyset) = 0$
- First probability axiom:

$$\text{pr}(E) \geq 0 \quad \text{for all } E \subseteq \Omega$$

- What is the largest possible probability?

- Assume that

$$\Omega = \{\text{red, black, green}\}$$

$$\text{then } \text{pr}(\{\text{red, black, green}\}) = 1$$

- Second probability axiom:

$$\text{pr}(\Omega) = 1$$

PROBABILITY AXIOMS – AXIOM III

- The third axiom covers the additivity of independent events

$$\text{pr}(E_1 \cup E_2 \cup \dots \cup E_n) = \sum_i^n \text{pr}(E_i)$$

if all E_i are independent

- Independence is not given if for example

$$E_1 = \{\text{black, green}\}, E_2 = \{\text{red, green}\}.$$

- In this case we have

$$\text{pr}(E_1 \cup E_2) \neq \text{pr}(E_1) + \text{pr}(E_2) > 1$$

Probability distribution (discrete case)

A probability distribution $\text{pr} : \mathbb{P}(\Omega) \rightarrow [0, 1]$ is a function that assigns a probability to each element of the powerset of Ω . In addition, it fulfills the probability axioms I-III, i.e.

- $\text{pr}(E) \geq 0$ for all $E \subseteq \Omega$
- $\text{pr}(\Omega) = 1$
- $\text{pr}(E_1 \cup E_2 \cup \dots \cup E_n) = \sum_i^n \text{pr}(E_i)$

where E_1, E_2, \dots, E_n are pairwise independent events.

COMPLEMENT AND SUM RULE

- There are several consequences of the probability axioms, one is the complement rule

$$\text{pr}(E^c) = \text{pr}(\Omega) - \text{pr}(E) = 1 - \text{pr}(E)$$

- For example, the probability of not observing *black* is given by

$$\text{pr}(\{\text{black}\}^c) = 1 - \text{pr}(\{\text{black}\})$$

- Another important consequence is the addition law or *sum rule*, given by

$$\text{pr}(A \cup B) = \text{pr}(A) + \text{pr}(B) - \text{pr}(A \cap B)$$

where $\text{pr}(A \cap B) = 0$ if A and B are independent

RANDOM VARIABLES

- Random variables (RVs) add another layer of formalism – Why do we need them?
- Assume we consider a more complex random experiment, i.e. we observe the roulette game for n rounds
- What is the probability of observing *black* in the i th round?
- To formalize this notion, we would associate a random variable X_i with the i th round and write

$$\text{pr}(X_i = \{\text{black}\})$$

- Similarly, we write

$$\text{pr}(X_i = \{\text{black}\}, X_j = \{\text{green}\})$$

for observing *black* in the i th round and *green* in the j round

- There exist two types of roulette wheels:
 - ▶ 1 green segment and 37 in total (what we considered)
 - ▶ 2 green segments and 38 in total
- Let X_1 correspond to the first type and X_2 the second
- We see that

$$\text{pr}(X_1 = \{\text{black}\}) \neq \text{pr}(X_2 = \{\text{black}\})$$

- **Random variables correspond to different types of distributions (probability assignments)**

Random Variable (RV)

A random variable $X : \Omega \rightarrow \mathbb{R}$ is a mapping from the sample space Ω to a measurable space, typically the real numbers \mathbb{R} . We denote by $\{X = x\}$ the *event* that X takes the value x and with $X \sim D$ that X has distribution D .

- Our previous notation, e.g. $X_1 = \{\text{black}\}$, is not correct
- We stick to this notation for simplicity
- If possible we avoid random variables, to simplify notation, i.e. we write

$$\text{pr}(\{\text{black}\}) = \text{pr}(X_1 = \{\text{black}\})$$

if unambiguous

CONDITIONAL PROBABILITIES

- A conditional probability is the probability of an event given that another event has happened or is known
- Assume we know that the ball has landed on a red segment. What is the probability that the segment has an even number?
- This is a conditional probability denoted as

$$\text{pr}(\{\text{even}\} | \{\text{red}\}) = ?$$

CONDITIONAL PROBABILITIES - DERIVATION

- Consider the following ingredients:
 - ▶ $\text{pr}(\{\text{even}\} | \{\text{red}\})$: The probability of an *even* segment, given that the ball has landed on a red segment
 - ▶ $\text{pr}(\{\text{red}\})$: The probability that we observe *red*
- What is the probability of *red* **and** of an *even* segment, given that the ball has landed on a red segment?

$$\text{pr}(\{\text{red}\})\text{pr}(\{\text{even}\} | \{\text{red}\})$$

- Logically, this is equivalent to asking: What is the probability of observing a *red* segment with an *even* number

$$\text{pr}(\{\text{red}\})\text{pr}(\{\text{even}\} | \{\text{red}\}) = \text{pr}(\{\text{even and red}\})$$

CONDITIONAL PROBABILITIES - DERIVATION

- Let A denote the set of all *red* segments
- Let B denote the set of all *even* segments
- The set of *red* segments with an *even* number is $A \cap B$
- Hence, we can rewrite

$$\text{pr}(\{\text{red}\})\text{pr}(\{\text{even}\} | \{\text{red}\}) = \text{pr}(\{\text{even and red}\})$$

as follows:

$$\begin{aligned}\text{pr}(A)\text{pr}(B | A) &= \text{pr}(A \cap (A^c \cup B)) \\ &= \text{pr}(A \cap B)\end{aligned}$$

- Note that this is equivalent to logic calculus:

$$\begin{aligned}A \wedge (A \rightarrow B) &= A \wedge (\neg A \vee B) \\ &= A \wedge B\end{aligned}$$

Conditional probability and Bayes theorem

The *conditional probability* of A given B is defined through

$$\text{pr}(A | B)\text{pr}(B) = \text{pr}(A \cap B) = \text{pr}(B | A)\text{pr}(A)$$

If $\text{pr}(B) > 0$ it follows that

$$\text{pr}(A | B) = \frac{\text{pr}(A \cap B)}{\text{pr}(B)} = \frac{\text{pr}(B | A)\text{pr}(A)}{\text{pr}(B)}$$

This is called *Bayes theorem*, where we call

- $\text{pr}(A | B)$: the posterior
- $\text{pr}(B | A)$: the likelihood
- $\text{pr}(A)$: the prior probability
- $\text{pr}(B)$: the evidence or marginal likelihood

Independence

Two events A and B are called *independent*, if

$$\text{pr}(A \cap B) = \text{pr}(A | B)\text{pr}(B) = \text{pr}(A)\text{pr}(B).$$

We also denote independence as $A \perp B$.

■ Example:

- ▶ We observe two rounds of roulette, associated with random variables X_1 , and X_2
- ▶ The probability of observing first *black* and then *red* is

$$\begin{aligned}\text{pr}(X_1 = \{\text{black}\}, X_2 = \{\text{red}\}) \\ &= \text{pr}(X_2 = \{\text{red}\} | X_1 = \{\text{black}\})\text{pr}(X_1 = \{\text{black}\}) \\ &= \text{pr}(X_2 = \{\text{red}\})\text{pr}(X_1 = \{\text{black}\})\end{aligned}$$

THE GAMBLER'S FALLACY

The Gambler's Fallacy [Hacking, 2001]

Consider our roulette game. The Fallacious Gambler reasons as follows:

- The roulette wheel is fair
- I have just observed 12 black spins in a row
- Since the wheel is fair, black and red come up equally often
- Hence, red has to come up pretty soon, I'd better start betting red

The gambler thinks that a sequence of twelve blacks makes it more likely that the wheel will step at red next time. If so, a past sequence affects future outcomes and the wheel is not fair. So trials would not be independent and the gambler's premises are inconsistent.

Law of total probability

Let B_1, B_2, \dots, B_n denote n mutually independent and exhaustive events, then

$$\text{pr}(A) = \sum_{i=1}^n \text{pr}(A \cap B_i) = \sum_{i=1}^n \text{pr}(A | B_i) \text{pr}(B_i)$$

- Example:

$$\text{pr}(\{\text{black}\}) = \text{pr}(\{\text{black and even}\}) + \text{pr}(\{\text{black and odd}\})$$

- The set of even and odd segments does not overlap (independence) and covers the full sample space (exhaustive)

EXAMPLES

Inductive logic [Hacking, 2001]

Inductive logic is about risky arguments. It analyses inductive arguments using probability. A risky argument can be a very good one, and yet its conclusion can be false.

- Probability assignments reflect beliefs about events
- They may come from simple distributions or complex models, such as neural networks
- Bayes theorem is used to derive probabilistic **if-statements**:
 - ▶ If B has happened, what is the probability of A ?

Medical testing

After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a serious disease, and that the test is 99% accurate (i.e., the probability of testing positive given that you have the disease is 0.99, as is the probability of testing negative given that you don't have the disease). The good news is that this is a rare disease, striking only one in 10,000 people. What are the chances that you actually have the disease?

- Let I denote a random variable indicating infection, i.e. $I = \text{true}$ if you are infected
- Let T denote a random variable associated with the test result, i.e. $T = \text{true}$ if the test is positive
- We know that:
 - ▶ $\text{pr}(I = \text{true}) = 1/10,000$
 - ▶ $\text{pr}(T = \text{true} | I = \text{true}) = 0.99$
 - ▶ $\text{pr}(T = \text{false} | I = \text{false}) = 0.99$

$$\begin{aligned} & \text{pr}(I = \text{true} \mid T = \text{true}) \\ &= \frac{\text{pr}(T = \text{true} \mid I = \text{true})\text{pr}(I = \text{true})}{\sum_i \text{pr}(T = \text{true} \mid I = i)\text{pr}(I = i)} \\ &= \frac{0.99 \cdot 1/10,000}{0.99 \cdot 1/10,000 + (1 - 0.99) \cdot (1 - 1/10,000)} \\ &= \frac{0.000099}{0.000099 + 0.009999} \\ &\approx 0.0098 \end{aligned}$$

Hence, the chance of actually having the disease is less than 1%.

THE MONTY HALL PROBLEM

The Monty Hall problem

On a game show, a contestant is told the rules as follows: There are three doors, labeled 1, 2, 3. A single prize has been hidden behind one of them with equal probability. You get to select one door. Initially your chosen door will not be opened. Instead, the gameshow host will open one of the other two doors, and he will do so in such a way as not to reveal the prize. For example, if you first choose door 1, he will then open one of doors 2 and 3, and it is guaranteed that he will choose which one to open so that the prize will not be revealed. At this point, you will be given a fresh choice of door: you can either stick with your first choice, or you can switch to the other closed door. All the doors will then be opened, and you will receive whatever is behind your final choice of door.

The Monty Hall problem

Imagine that the contestant chooses door 1 first; then the gameshow host opens door 3, revealing nothing behind the door, as promised. Should the contestant

- stick with door 1,
- switch to door 2,
- does it make a difference?

THE MONTY HALL PROBLEM

- Let $P \in \{1, 2, 3\}$ denote a random variable associated with the location of the price
- Let $D \in \{1, 2, 3\}$ denote the door that has been opened by the gameshow host
- A priori we have

$$\text{pr}(P = i) = 1/3$$

- We are interested in the posterior probability

$$\text{pr}(P = i | D = 3) = \frac{\text{pr}(D = 3 | P = i)\text{pr}(P = i)}{\text{pr}(D = 3)}$$

THE MONTY HALL PROBLEM

■ The case $i = 1$

$$\begin{aligned}\text{pr}(P = 1 | D = 3) &= \frac{\text{pr}(D = 3 | P = 1)\text{pr}(P = 1)}{\sum_i \text{pr}(D = 3 | P = i)\text{pr}(P = i)} \\ &= \frac{1/2 \cdot 1/3}{1/2 \cdot 1/3 + 1 \cdot 1/3 + 0 \cdot 1/3} \\ &= 1/3\end{aligned}$$

■ The case $i = 2$

$$\begin{aligned}\text{pr}(P = 2 | D = 3) &= \frac{\text{pr}(D = 3 | P = 2)\text{pr}(P = 2)}{\sum_i \text{pr}(D = 3 | P = i)\text{pr}(P = i)} \\ &= \frac{1 \cdot 1/3}{1/2 \cdot 1/3 + 1 \cdot 1/3 + 0 \cdot 1/3} \\ &= 2/3\end{aligned}$$

THE MONTY HALL PROBLEM

- The case $i = 3$

$$\begin{aligned}\text{pr}(P = 1 | D = 3) &= \frac{\text{pr}(D = 3 | P = 1)\text{pr}(P = 1)}{\sum_i \text{pr}(D = 3 | P = i)\text{pr}(P = i)} \\ &= \frac{0 \cdot 1/3}{1/2 \cdot 1/3 + 1 \cdot 1/3 + 0 \cdot 1/3} \\ &= 0\end{aligned}$$

- Hence, we have the posterior distribution

$$\text{pr}(P = i | D = 3) = (1/3, 2/3, 0)$$

- Switching the door increases the probability of getting the price

PROBABILITY DISTRIBUTIONS

Bernoulli distribution

Let X be a random variable taking values in $\{0, 1\}$. If X follows a Bernoulli distribution with parameter p , i.e.

$$X \sim \text{Bernoulli}(p)$$

then

$$\text{pr}(X = 1) = p.$$

- Flipping a coin once can be modeled using a Bernoulli distribution
- The coin is fair if $p = 1/2$

Categorical or multinoulli distribution

Let X be a random variable taking values in $\{0, 1, \dots, k\}$ for any integer $k > 0$. If X follows a categorical distribution with parameters $p = (p_1, \dots, p_k)$ such that $\sum_i p_i = 1$, i.e.

$$X \sim \text{Categorical}(p)$$

then

$$\text{pr}(X = i) = p_i.$$

- The categorical distribution is the extension of the Bernoulli distribution to k outcomes

Geometric distribution

Let X be a random variable taking values in $\{1, 2, 3, \dots\}$. If X follows a geometric distribution with parameter $p \in [0, 1]$, i.e.

$$X \sim \text{Geometric}(p)$$

then

$$\text{pr}(X = k) = (1 - p)^{k-1}p.$$

- The probability distribution of the number X of Bernoulli trials needed to get one success
- It gives the probability that the first occurrence of success requires k independent trials, each with success probability p .

Normal or Gaussian distribution

Let X be a random variable taking values in \mathbb{R} . If X follows a normal distribution with parameters $\mu \in \mathbb{R}$ and $\sigma > 0$, i.e.

$$X \sim \text{Normal}(\mu, \sigma)$$

then

$$\text{pr}(X \in A) = \int_A f_{\mu, \sigma}(x) dx,$$

where $f_{\mu, \sigma}$ is the normal density function

$$f_{\mu, \sigma}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right).$$

- The probability distribution of continuous random variables is usually defined through density functions
- For continuous distributions we have $\text{pr}(X = x) = 0$ for all $x \in \mathbb{R}$, i.e. the probability that a single real value is observed is always zero

PARAMETER ESTIMATION

PARAMETER ESTIMATION

- Assume we observed n realisations $\mathbf{x} = (x_1, \dots, x_n)$ from a known distribution with unknown parameters θ
- How can we estimate the values of θ ?
- Bayes theorem

$$\text{pr}(\theta | \mathbf{x}) = \frac{\text{pr}(\mathbf{x} | \theta)\text{pr}(\theta)}{\text{pr}(\mathbf{x})}$$

- We take the value with highest probability, i.e.

$$\hat{\theta} = \arg \max_{\theta} \text{pr}(\theta | \mathbf{x})$$

this is called the *maximum a-posteriori (MAP) estimate*

- Constants can be dropped when computing the MAP, i.e.

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \text{pr}(\theta | \mathbf{x}) \\ &= \arg \max_{\theta} \frac{\text{pr}(\mathbf{x} | \theta) \text{pr}(\theta)}{\text{pr}(\mathbf{x})} \\ &= \arg \max_{\theta} \text{pr}(\mathbf{x} | \theta) \text{pr}(\theta) \\ &= \arg \max_{\theta} [\log \text{pr}(\mathbf{x} | \theta) + \log \text{pr}(\theta)]\end{aligned}$$

since $\text{pr}(\mathbf{x})$ does not depend on θ . We can apply the logarithm, because it is a monotonic (*order preserving*) function, which does not change the position of the maximum

- If we assume that we have no prior information on θ , i.e. $\text{pr}(\theta)$ is uniform (constant), then we obtain the *maximum likelihood (ML) estimate*

$$\hat{\theta} = \arg \max_{\theta} \text{pr}(\mathbf{x} | \theta)$$

PARAMETER ESTIMATION – CONTINUOUS VARIABLES

- For continuous variables we know that

$$\text{pr}(x | \theta) = 0$$

and therefore also $\text{pr}(x)$ is zero, which causes the posterior distribution to be undefined

- There exists a Bayes theorem for densities

$$f(\theta | x) = \frac{f(x | \theta)f(\theta)}{f(x)}$$

- Hence, the MAP for continuous variables is simply

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} f(x | \theta)f(\theta) \\ &= \arg \max_{\theta} [\log f(x | \theta) + \log f(\theta)]\end{aligned}$$

NORMAL DISTRIBUTION

- Assume we observed n realisations $x = (x_1, \dots, x_n)$ from $X \sim \text{Normal}(\mu, \sigma)$ with unknown μ and σ
- Furthermore, let's assume we have no prior information about μ and σ , i.e. the prior probability $\text{pr}(\mu, \sigma)$ is uniform
- We derive the ML estimate

$$\begin{aligned}\hat{\mu}, \hat{\sigma} &= \arg \max_{\mu, \sigma} f_{\mu, \sigma}(x) \\ &= \arg \max_{\mu, \sigma} \prod_{i=1}^n f_{\mu, \sigma}(x_i) \\ &= \arg \max_{\mu, \sigma} \sum_{i=1}^n \log f_{\mu, \sigma}(x_i)\end{aligned}$$

- We derive the ML estimate

$$\begin{aligned}\hat{\mu}, \hat{\sigma} &= \arg \max_{\mu, \sigma} \sum_{i=1}^n \log f_{\mu, \sigma}(x_i) \\ &= \arg \max_{\mu, \sigma} + \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \\ &= \arg \max_{\mu, \sigma} -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\end{aligned}$$

NORMAL DISTRIBUTION

- We derive the ML estimate of μ

$$\frac{\partial}{\partial \mu} \left[-\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right] = 0$$

$$\Rightarrow -\frac{\partial}{\partial \mu} \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\Rightarrow \frac{\partial}{\partial \mu} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\Rightarrow \sum_{i=1}^n -2(x_i - \mu) = 0$$

$$\Rightarrow \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

NORMAL DISTRIBUTION

- We derive the ML estimate of σ^2

$$\frac{\partial}{\partial \sigma^2} \left[-\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right] = 0$$

$$\Rightarrow -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\Rightarrow \frac{1}{2\sigma^2} \left[-n + \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right] = 0$$

$$\Rightarrow \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 = n$$

$$\Rightarrow \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

- [Hacking, 2001]

REFERENCES



HACKING, I. (2001).

AN INTRODUCTION TO PROBABILITY AND INDUCTIVE LOGIC.

Cambridge university press.



KOLMOGOROFF, A. (1933).

GRUNDBEGRIFFE DER WAHRSCHEINLICHKEITSRECHNUNG.